

# Faster coreset construction for projective clustering *via* low rank approximation

Rameshwar Pratap<sup>1</sup> and Sandeep Sen<sup>2</sup>

<sup>1</sup> IIIT Bangalore, India

rameshwar.pratap@gmail.com

<sup>2</sup> IIT Delhi, India

ssen@cse.iitd.ernet.in

**Abstract.** In this work, we present randomized coreset construction for projective clustering that involves computing a set of  $k$  closest  $j$ -dimensional linear (affine) subspaces of a given set of  $n$  vectors in  $d$  dimensions. Let  $A \in \mathbb{R}^{n \times d}$  be an input matrix, then an earlier deterministic coreset construction of Feldman *et. al* [10] relied on computing the SVD of  $A$ . The best known algorithms for SVD require  $nd \min\{n, d\}$  time, which may not be feasible for large values of  $n$  and  $d$ .

We present a coreset construction by projecting the matrix  $A$  on some orthonormal vectors that closely approximate the right singular vectors of  $A$ . As a consequence, when the values of  $k$  and  $j$  are small, we are able to achieve a significantly faster algorithm as compared to Feldman *et. al* [10], while maintaining almost the same approximation. We also benefit in terms of space as well as exploit the sparsity of the input dataset. Another advantage of our approach is that it can be constructed in a streaming setting quite efficiently.

## 1 Introduction

*Succinct representation of Big data – Coreset:* Recent years have witnessed a dramatic increase in our ability to collect data from various sources. This data flood has surpassed our ability to understand, analyse and process them. *Big data* is a new terminology that has become quite popular to identify such datasets that are difficult to analyse with the current available technologies. One possible approach to manage such large volume of datasets is to keep a succinct summary of the datasets such that it approximately preserves the required properties of the original datasets. This notion was initially formalised by Agrawal *et al.* [1] in the context of approximating various descriptors of the extent of a point set. They derived  $\epsilon$ -approximation algorithms for computing smallest enclosing ball, computing diameter, width etc., and coined the term *coreset* for such summaries of the input.

Intuitively, a coreset can be considered as a semantic compression of the input. More precisely, a coreset is a weighted subset of the data such that the quality of any clustering evaluated on the coreset closely approximates the quality on the full data set. Consider a set  $Q$  (possibly of infinite size) of query

shapes (for example: subspaces, set of points, set of lines etc.), then for every shape  $q \in Q$ , the sum of distances from  $q$  to the input points, and the sum of distances from  $q$  to the points in the coreset, is approximately the same. If the query set belongs to some particular candidate query set, then such coreset is called as *weak* coreset (see a survey on weak coreset [17]); and if the coreset approximates the distances from all possible (potentially infinite) query shapes, then it is called as *strong* coreset.

Coresets are practical and flexible tool which require no or very minimal assumption on the data. Although the analysis and proof techniques for coreset construction are a bit involved, and require tools from geometry and linear algebra, the resulting coreset construction algorithms are very easy to implement. Another important property of coresets is that they can be constructed in streaming and distributed setting quite efficiently. This is due to the fact that unions of coresets are coresets, and coresets of coresets are coreset [12]. Also, using these properties it is possible to construct coresets in a tree-wise fashion which can be parallelized in a Map-Reduce style [10].

Coreset constructions have been studied extensively for various data analysis task. Their construction and analysis techniques mostly include geometric approximations and linear algebraic approach. Geometric coreset constructions are usually done in an iterative fashion. To start with, some weights are assigned to each point, and in each iteration their weights are refined as per the significance of the points with respect to the problem. Linear algebraic approach of coreset construction usually done in two steps – dimensionality reduction, and cardinality reduction. The dimension reduction step of the coreset construction includes projecting points in a low dimension space such that the original geometry of points is also preserved in a low dimensional space. These projection techniques includes SVD decomposition, random projections, row/column subset selections, or any combinations of these (see [10,5]). The cardinality reduction step includes contracting the input size *via* sampling or other geometric analysis approach on the reduced dimension instance of the input. We refer readers to a recent survey article of Jeff M. Phillips [18] on coreset and sketches covering the above, and an earlier survey by Agarwal *et al.* [2] on geometric coreset construction.

In this work, we focus on the dimension reduction step of coreset construction for projective clustering problem. In the paragraph below, we discuss the motivation behind the projective clustering problem.

*Projective clustering:* Clustering is one of the most popular technique for analyzing large data, and is widely used in many areas such as classification, unsupervised learning, data mining, indexing, pattern recognition. Many popular clustering algorithm such as  $k$ -mean [16], BIRCH [21], DBSCAN [6] are full dimensional – they give equal importance to all the dimensions while computing the distance between two points. These clustering algorithms works well in low dimensional datasets, however, due to the “*curse of dimensionality*” scale poorly in high dimension. Moreover, in high dimensional datasets a full dimensional distance might not be appropriate as farthest neighbour of a point is expected to be roughly as close as its nearest neighbour [14]. These problems are often

handled *via* methods such as Principal component analysis (PCA) or Johnson-Lindenstrauss lemma [15] by finding a low dimensional representation of the data obtained by projecting all points on a subspace so that the information loss is minimized. However, projecting all the points in a single low dimensional subspace may not be appropriate when different clusters lie in different subspaces. This motivates the study of projective clustering which involves finding clusters along different subspaces. Projective clustering algorithms have been widely applicable for indexing [4] and pattern discovery [3] in high dimensional datasets.

### 1.1 Our contribution

With the above motivation we study the dimension reduction step of coresset construction for projective clustering problem. We first briefly describe the subspace and projective clustering problems. In a  $j$ -subspace clustering problem, given a set of  $n$   $d$  dimensional vectors, denoted by  $A \in \mathbb{R}^{n \times d}$ , the problem is to find a  $j$ -dimensional subspace such that it minimizes the sum of squared distances over every  $j$ -dimensional subspace. Further, in the problem of linear (affine)  $(k, j)$ -projective clustering, the goal is to find a closed set  $\mathcal{C}$  which is the union of  $k$  linear (affine) subspaces each of dimension  $j$ , such that it minimizes the sum of squared distances over every possible choice of  $\mathcal{C}$  (see Definitions 8,9).

Feldman *et al.* [10] presented a deterministic coresset construction for these clustering problems. Their coresset construction relies on projecting the rows of  $A$  on the first few right singular values of  $A$ . However, the main drawback of their construction is that it requires computing the SVD of  $A$  which is expensive for large values of  $n$  and  $d$ . Recently, Cohen *et al.* [5] suggested “*projection-cost-preserving-sketch*” for various clustering problem. Their sketches are essentially the dimensionality reduction step of the coresset construction. Using a low rank approximation of  $A$ , they suggested a faster coresset construction for the subspace clustering problem. However, it was not clear that how their techniques can be extended for projective clustering problem<sup>3</sup>. In this work, we extend their techniques and obtain a faster dimension reduction for projective clustering, and as a consequence, a faster coresset construction for projective clustering problem. In Section 3, we first revisit the techniques for subspace clustering problem, and in Section 4 we present our coresset construction for projective clustering problem. We state our main result as follows: (In the following theorem,  $\mathbf{nnz}(A)$  denotes the number of non-zero entries of  $A$ .)

**Theorem 1.** *Let  $A \in \mathbb{R}^{n \times d}$ ,  $\epsilon \in (0, 1)$ , and  $j, k$  be two integers less than  $(d-1)$ , and  $(n-1)$  respectively such that  $k(j+1) \leq d-1$ . Then there is a randomized algorithm which outputs a matrix  $A^*$  of rank  $O\left(\frac{k(j+1)}{\epsilon^2}\right)$  such that for every non-empty closed set  $\mathcal{C}$ , which is the union of  $k$  linear (affine) subspaces each of dimension at most  $j$ , the following holds w.h.p.*

$$|(\text{dist}^2(A^*, \mathcal{C}) + \Delta^*) - \text{dist}^2(A, \mathcal{C})| \leq \epsilon \text{dist}^2(A, \mathcal{C}).$$

---

<sup>3</sup> private communication

Where,  $j^* = k(j + 1)$ ;  $\Delta^* = \|A - A^{O(\frac{j^*}{\epsilon^2})}\|_F^2$ ;  $\text{dist}^2(A, C)$  denotes the sum of squared distances from each row of  $A$  to its closest point in  $C$ ; and  $A^{O(\frac{j^*}{\epsilon^2})}$  is the best rank  $O(\frac{j^*}{\epsilon^2})$  approximation of  $A$ . The expected running time of the algorithm is  $\tilde{O}\left(\text{nnz}(A)\frac{j^*}{\epsilon^3} + (n + d)\frac{j^{*2}}{\epsilon^6} + \frac{ndj^*}{\epsilon^2}\right)$ .<sup>4</sup>

*Remark 1.* We develop our coresets by projecting points on some orthonormal vectors that closely approximate the right singular vectors of  $A$ , and we obtain them using the algorithm of Sarlós [19]. The expected running time of our algorithm is better than the corresponding deterministic algorithm of [10] when  $n \geq d$  and  $j^* = o(n)$ , or, when  $n < d$  and  $j^* = o(d)$ , where  $j^* = k(j + 1)$ . Further, as the coreset construction time depends on number of non-zero entries of the matrix, our algorithm is substantially faster for sparse data matrices.

*Remark 2.* An advantage of our coresets is that it can be constructed in the pass efficient streaming model [13] - where access to the input is limited to only a constant number of sequential passes. We construct our coreset by projecting the matrix  $A$  on orthonormal vectors, that closely approximate the right singular vectors of  $A$ , our algorithm requires only two passes over the data in order to compute those orthonormal vectors using [19].

## 1.2 Related work

Coreset construction has been studied extensively for the problem of  $j$ -subspace clustering. However, we will discuss a few of them that are more relevant to our work. Feldman *et al.* [7] developed a strong coreset whose size is exponential in  $d, j$ , logarithmic in  $n$ , and their coreset construction requires  $O(n)$  time. Feldman *et al.* [9] improved their earlier result [7] and developed a coreset of size logarithmic in  $n$ , linear in  $d$ , and exponential in  $j$ . However, the construction requires  $O(ndj)$  time. In [8] Feldman and Langberg showed a coreset construction of size polynomial in  $j$  and  $d$  (independent of  $n$ ). Feldman *et al.* [10] presented a novel coreset construction for subspace and projective clustering. They showed that the sum of square Euclidean distance from  $n$  rows of  $A \in \mathbb{R}^{n \times d}$  to any  $j$ -dimensional subspace can be approximated upto  $(1 + \epsilon)$  factor, with an additive constant which is the sum of a few last singular values of  $A$ , by projecting the points on the first  $O(j/\epsilon)$  right singular vectors of  $A$ . Thus, they able to show the dimension reduction from  $d$  to  $O(j/\epsilon)$ . They also showed  $O(k(j + 1)/\epsilon^2)$  dimension reduction for  $(k, j)$ -projective clustering problem. Recently, for  $j$ -subspace clustering, Cohen *et al.* [5] improved the construction of [10] by using only first  $\lceil j/\epsilon \rceil$  right singular vectors, which is an improvement over [10] by a constant factor.

Sariel Har-Peled [11] showed that for projective clustering problem it is not possible to get a strong coreset of size sublinear in  $n$  - even for family of pair of planes in  $\mathbb{R}^3$ . However, Varadarajan *et al.* [20] showed a sublinear size coreset for projective clustering on a restricted setting - when points are on an integer

<sup>4</sup> Here,  $\tilde{O}$  is the asymptotic notation that ignores logarithmic factors.

grid, and the largest coordinate of any point is bounded by a polynomial in  $n$  and  $d$ .

### 1.3 Organization of the paper

In Section 2, we present the necessary notations, definition and linear algebra background are used in the various proofs in the paper. In Section 3, we revisit the result of [5] and discuss the coresset construction for subspace clustering using their techniques. In Section 4, we extend the result of Section 3 and we present the coresset construction for projective clustering problem. We conclude our discussion and state some open questions in Section 5.

## 2 Preliminaries

Notations	
$A = U\Sigma V^T$	columns of $U, V$ are orthonormal and called as left and right singular vectors of $A$ ; $[\Sigma]$ is a diagonal matrix having the corresponding singular values
$A^{(m)} = U\Sigma^{(m)}V^T$	$\Sigma^{(m)}$ is the diagonal having the $m$ largest entries of $\Sigma$ , and 0 otherwise
$[X]_{d \times j}$	$j$ orthonormal columns represent a $j$ -dimensional subspace $L$ in $\mathbb{R}^d$
$[X^\perp]_{d \times (d-j)}$	$(d-j)$ dimensional subspace $L^\perp$ orthogonal to subspace $L$
$\pi_{\mathcal{S}}(A)$	matrix formed by projecting $A$ on the row span of $\mathcal{S}$
$\pi_{\mathcal{S},k}(A)$	the best rank- $k$ approximation of $A$ with its rows projected on the row span of $\mathcal{S}$
$A^{(k)}$	the best rank- $k$ approximation of $A$
$\text{nnz}(A)$	the number of non-zero entries of $A$

Below we present some necessary linear algebra background. We first present some basic properties of Frobenius norm of a matrix. Then we define SVD (singular value decomposition) of a matrix, and its basic properties. Then, we discuss the expression about the distance of a point, and sum of square distances of the rows of matrix - from a subspace and a closed set.

**Fact 1 (Frobenius norm and its properties)** *Let  $A \in \mathbb{R}^{n \times d}$ , then square of Frobenius norm of  $A$  is defined as the sum of the absolute squares of its elements, i.e.  $\|A\|_F^2 = \sum_{i=1}^n \sum_{j=1}^d a_{i,j}^2$ . Further, if  $\{\sigma_i\}_{i=1}^d$  are singular values of  $A$ , then  $\|A\|_F^2 = \sum_{i=1}^d \sigma_i^2$ . Also, if  $\text{tr}(A)$  be the trace of the matrix  $A$  then  $\|A\|_F^2 = \text{tr}(A^T A)$ .*

**Fact 2** *Let  $AX$  be the projection of points of  $A$  on the  $j$ -dimensional subspace  $L$  represented by a matrix  $X$ . We can also write the projection of the points in the rows of  $A$  to  $L$  as  $AXX^T$ , these projected points are still  $d$ -dimensional, but lie within the  $j$ -dimensional subspace. Further,  $\|AX\|_F^2 = \|AXX^T\|_F^2$ .*

*The Singular Value Decomposition:* A matrix  $A \in \mathbb{R}^{n \times d}$  of rank at most  $r$  can be written due to its SVD decomposition as  $A = \sum_{i=1}^r \sigma_i u^{(i)} v^{(i)T}$ . Here,  $u^{(i)}$  and  $v^{(i)}$  are  $i$ -th orthonormal columns of  $U$  and  $V$  respectively, and  $\sigma_1 \geq \sigma_2, \dots, \sigma_r \geq 0$ . Also,  $u^{(i)T} A = \sigma_i v^{(i)T}$ , and  $A v^{(i)} = \sigma_i u^{(i)}$  for  $1 \leq i \leq r$ . Further, the matrix  $A^{(k)}$  that minimizes  $\|A - B\|_F$  among all matrices  $B$  (of rank at most  $k$ ) is given by  $A^{(k)} = \sum_{i=1}^k A v^{(i)} v^{(i)T}$  - i.e. by projecting  $A$  on the first  $k$  right singular vectors of  $A$ .

*$l_2$  distances to a subspace:* Let  $L$  be a  $j$ -dimensional subspace in  $\mathbb{R}^d$  represented by an orthonormal matrix  $X \in \mathbb{R}^{d \times j}$ . Then, for a point  $p \in \mathbb{R}^d$ ,  $\|p^T X\|_F^2$  is the square of the length of projection of the point  $p$  on the subspace  $L$ . Similarly, given a matrix  $A \in \mathbb{R}^{n \times d}$ ,  $\|AX\|_F^2$  is the sum of square of the length of projection of the points (rows) of  $A$  on the subspace  $L$ . Let  $L^\perp$  be the orthogonal complement of  $L$  represented by an orthonormal matrix  $X^\perp \in \mathbb{R}^{d \times (d-j)}$ . Then,  $\|AX^\perp\|_F^2$  is the sum of square of distances of the points of  $A$  from  $L$ .

*$l_2$  distance to a closed set:* Let  $S \in \mathbb{R}^d$  be a closed set and  $p$  be a point in  $\mathbb{R}^d$ . We define the  $l_2$  distance between  $p$  and  $S$  by  $\text{dist}^2(p, S) := \min_{s \in S} \text{dist}^2(p, s)$ , i.e., the smallest distance between  $p$  and any element  $s \in S$ . If  $S$  consists of union of  $k$ ,  $j$ -dimensional subspaces  $L_1, \dots, L_k$ , then  $\text{dist}^2(p, S)$  denotes the distance from  $p$  to the closest set  $S$ . Similarly, given a matrix  $A \in \mathbb{R}^{n \times d}$ ,  $\text{dist}^2(A, S) := \sum_{i=1}^n \text{dist}^2(A_{i*}, S)$ . Here,  $A_{i*}$  denotes the  $i$ th row of  $A$ .

*Pythagorean theorem:* Let  $A \in \mathbb{R}^{n \times d}$ ,  $L$  be a  $j$ -dimensional subspace in  $\mathbb{R}^d$  represented by an orthonormal matrix  $X \in \mathbb{R}^{d \times j}$ , and  $L^\perp$  be the orthogonal complement of the subspace  $L$  represented by an orthonormal matrix  $X^\perp \in \mathbb{R}^{d \times (d-j)}$ . Then by Pythagorean theorem we have  $\|A\|_F^2 = \|AX\|_F^2 + \|AX^\perp\|_F^2$ . Further, if  $\mathcal{C}$  is a closed set, then due to the Pythagoras theorem we have  $\text{dist}^2(A, \mathcal{C}) = \|AX^\perp\|_F^2 + \text{dist}^2(AXX^T, \mathcal{C})$ . We will use the following fact in our analysis which hold true due to Pythagorean theorem.

**Fact 3** Let  $A \in \mathbb{R}^{n \times d}$ , and  $X \in \mathbb{R}^{d \times j}$  be a matrix having orthonormal columns, then due to the Pythagorean theorem we have  $\|A - AXX^T\|_F^2 = \|A\|_F^2 - \|AXX^T\|_F^2$ . Similarly, it can be shown that

$$(A - AXX^T)^T (A - AXX^T) = A^T A - (AXX^T)^T (AXX^T).$$

In the following, we state some facts from elementary linear algebra which are required for deriving the correctness of our result.

**Fact 4** For a square matrix  $M \in \mathbb{R}^{n \times n}$ ,  $\text{tr}(M)$  is the sum of all its diagonal entries. Further, for matrices  $A \in \mathbb{R}^{n \times d}$ ,  $B \in \mathbb{R}^{d \times n}$  due to the cyclic property of the  $\text{tr}$  function, we have  $\text{tr}(AB) = \text{tr}(BA)$ . Also for square matrices  $M, N \in \mathbb{R}^{n \times n}$ , due to the linear property of the  $\text{tr}$  function, we have  $\text{tr}(M \pm N) = \text{tr}(M) \pm \text{tr}(N)$ .

**Fact 5** A symmetric matrix  $M \in \mathbb{R}^{n \times n}$  is positive semidefinite if  $x^T M x > 0$  for all  $x \in \mathbb{R}^n$ . A matrix  $M$  is positive semidefinite then the following two statements are equivalent:

- there is a real nonsingular matrix  $N$  such that  $M = N^T N$ ,
- all eigenvalues of  $M$  are nonnegative.

**Fact 6** Let  $A \in \mathbb{R}^{n \times d}$  and  $U \Sigma V^T$  be the SVD of  $A$ . Then, the first  $j$  columns of  $V$  span a subspace that minimizes the sum of square distances of the vectors in  $A$  from all  $j$ -dimensional subspace, and this sum is  $\sum_{i=j+1}^d \sigma_i^2$ . Thus, for any  $j$ -dimensional subspace represented by an orthonormal matrix  $X$ , we have  $\|AX^\perp\|_F^2 \geq \sum_{i=j+1}^d \sigma_i^2$ .

**Fact 7** Let  $M \in \mathbb{R}^{d \times l}$ , and  $X \in \mathbb{R}^{d \times k}$  be an orthonormal matrix. Then, by elementary linear algebra we have,  $\|XX^T M\|_F^2 \leq \|M\|_F^2$ .

In the following, we state the definitions of subspace and projective clustering.

**Definition 8 (Subspace clustering)** Let  $A \in \mathbb{R}^{n \times d}$  and  $j$  be an integer less than  $d$ . Then, the problem of  $j$ -subspace clustering is to find a  $j$ -dimensional subspace  $L$  of  $\mathbb{R}^d$  that minimizes the  $\text{dist}^2(A, L)$ . In other words, the goal is to find a matrix  $X^\perp \in \mathbb{R}^{d \times (d-j)}$  having orthonormal columns that minimizes  $\|AX^\perp\|_F^2$  over every such possible matrix  $X^\perp$ .

**Definition 9 (linear (affine)  $(k, j)$ -projective clustering)** Let  $A \in \mathbb{R}^{n \times d}$ ,  $j$  be an integer less than  $d$ , and  $k$  be an integer less than  $n$ . Then, the problem of linear (affine)  $(k, j)$ -projective clustering is to find a closed set  $\mathcal{C}$ , which is the union of  $k$  linear (affine) subspaces  $\{L_1, \dots, L_k\}$  each of dimension at most  $j$ , such that it minimizes the  $\text{dist}^2(A, \mathcal{C})$ , over every possible choice of  $\mathcal{C}$ .

**Theorem 2 (Low-rank approximation by [19]).** Let  $A \in \mathbb{R}^{n \times d}$ , and  $\pi(\cdot)$  denote the projection operators stated in the notation table. If  $\epsilon \in (0, 1]$  and  $\mathcal{S}$  is an  $(r \times n)$  Johnson-Lindenstrauss matrix with i.i.d. zero-mean  $\pm 1$  entries and  $r = O\left(\left(\frac{m}{\epsilon} + m \log m\right) \log \frac{1}{\delta}\right)$ , then with probability at least  $1 - \delta$  it holds that

$$\|A - \pi_{\mathcal{S}A, m}(A)\|_F^2 \leq (1 + \epsilon) \|A - A^{(m)}\|_F^2.$$

Further, computing the singular vectors spanning  $\pi_{\mathcal{S}A, m}(A)$  in two passes<sup>5</sup> over the data requires  $O(\text{nnz}(A)r + (n + d)r^2)$  time.

For our analysis, we will use a weak triangle inequality which can be stated as follows:

**Lemma 10 (Lemma 7.1 of [10])** For any  $\epsilon \in (0, 1)$ , a closed set  $\mathcal{C}$ , and two points  $p, q \in \mathbb{R}^d$ , we have

$$|\text{dist}^2(p, \mathcal{C}) - \text{dist}^2(q, \mathcal{C})| \leq \frac{12\|p - q\|^2}{\epsilon} + \frac{\epsilon}{2} \text{dist}^2(p, \mathcal{C}).$$

<sup>5</sup> Two passes are required as we first multiply  $A$  on the right with a Johnson-Lindenstrauss matrix  $\mathcal{S}$ , and then we project the rows of  $A$  again onto the row span of  $\mathcal{S}A$ .

### 3 Faster coreset construction for subspace clustering

In this section by revisiting the results of Cohen *et al.* [5], we present a randomized coreset construction for subspace clustering. The deterministic coreset construction of Feldman *et al.* [10] for subspace clustering problem relies on projecting the input matrix on its first few right singular vectors – projecting the rows of  $A$  on first few right singular vectors of  $A$  – which requires SVD computation of  $A$ . Cohen *et al.* [5] suggested that projecting the rows of  $A$  some orthonormal vectors that closely approximate the right singular vectors of  $A$  (obtain via e.g. [19]) does also satisfies the required properties of coreset *w.h.p.*, and as a consequence, gives a faster construction.

**Theorem 3 (Adapted from Theorem 8 of [5]).** *Let  $X \in \mathbb{R}^{d \times j}$  be an orthonormal matrix representing a subspace  $L$ , let  $X^\perp \in \mathbb{R}^{d \times (d-j)}$  be the orthonormal matrix representing the orthogonal complement of  $L$ ,  $\epsilon \in (0, 1)$ ,  $\delta \in (0, 1)$ ,  $m = \lceil \frac{j}{\epsilon} \rceil$ ,  $\Delta = \|A - A^{(m)}\|_F^2$ , and  $\tilde{A}$  is a rank  $m$  approximation of  $A$  satisfying Theorem 2. Then, the following is true with probability at least  $1 - \delta$ :*

$$0 \leq \left| \|\tilde{A}X^\perp\|_F^2 + \Delta - \|AX^\perp\|_F^2 \right| \leq 2\epsilon \|AX^\perp\|_F^2.$$

*Proof.* Using a result of Sarlós [19], we get a rank  $m$  approximation of  $A$ . If  $\mathcal{S}$  is an  $(r \times n)$  JL matrix, where  $r = O\left(\left(\frac{m}{\epsilon} + m \log m\right) \log \frac{1}{\delta}\right)$  (as stated in preliminaries, see Theorem 2). then the following holds true with probability at least  $1 - \delta$ :

$$\|A - \pi_{\mathcal{S}A, m}(A)\|_F^2 \leq (1 + \epsilon) \|A - A^{(m)}\|_F^2. \quad (1)$$

Here,  $A^{(m)}$  is the best  $m$  rank approximation of  $A$ . Let  $R'$  be the matrix having the first  $m$  right singular vectors of  $\pi_{\mathcal{S}A}(A)$ , and let we denote  $AR'R^T$  by  $\tilde{A}$ , then by Equation 1, the following holds true with probability at least  $1 - \delta$ :

$$\|A - \tilde{A}\|_F^2 \leq (1 + \epsilon) \|A - A^{(m)}\|_F^2 \quad (2)$$

In the following we show an upper bound on the following expression:

$$\begin{aligned} & \left| \|\tilde{A}X^\perp\|_F^2 + \Delta - \|AX^\perp\|_F^2 \right| \\ &= \left| \|\tilde{A}\|_F^2 - \|\tilde{A}X\|_F^2 + \|A - A^{(m)}\|_F^2 - \|A\|_F^2 + \|AX\|_F^2 \right| \end{aligned} \quad (3)$$

$$= \left| \|\tilde{A}\|_F^2 - \|\tilde{A}X\|_F^2 + \|A\|_F^2 - \|A^{(m)}\|_F^2 - \|A\|_F^2 + \|AX\|_F^2 \right| \quad (4)$$

$$\begin{aligned} &= \left| \|\tilde{A}\|_F^2 - \|A^{(m)}\|_F^2 - \|\tilde{A}X\|_F^2 + \|AX\|_F^2 \right| \\ &\leq \left| \|A^{(m)}\|_F^2 - \|A^{(m)}\|_F^2 + \|AX\|_F^2 - \|\tilde{A}X\|_F^2 \right| \\ &= \left| \|AX\|_F^2 - \|\tilde{A}X\|_F^2 \right| \leq 2\epsilon \|AX^\perp\|_F^2 \end{aligned} \quad (5)$$

Equality 3 follows from pythagoras theorem; Equality 4 follows from Fact 3, where  $A^{(m)} = AV'V'^T$ , and  $V' \in \mathbb{R}^{d \times m}$  having  $m$  columns from the first  $m$  right singular vectors of  $A$ ; finally Inequality 5 holds from Lemma 11.



**Lemma 11 (Adapted from Lemma 5 of [5])** *Let  $A \in \mathbb{R}^{n \times d}$ ,  $\tilde{A}$  is a rank  $m$  approximation of  $A$  satisfying Equation 2, then  $0 \leq \|AX\|_F^2 - \|\tilde{A}X\|_F^2 \leq 2\epsilon\|AX^\perp\|_F^2$ .*

*Proof.* We first express the term  $\|AX\|_F^2 - \|\tilde{A}X\|_F^2$  in the form of tr function:

$$\begin{aligned} \|AX\|_F^2 - \|\tilde{A}X\|_F^2 &= \text{tr}((AX)^T(AX)) - \text{tr}((\tilde{A}X)^T\tilde{A}X) \\ &= \text{tr}(X^T A^T A X) - \text{tr}(X^T \tilde{A}^T \tilde{A} X) \\ &= \text{tr}(X^T (A^T A - \tilde{A}^T \tilde{A}) X) = \text{tr}(X X^T (A^T A - \tilde{A}^T \tilde{A})). \end{aligned} \quad (6)$$

The above equalities follows due to definition of tr function -  $\|A\|_F^2 = \text{tr}(A^T A)$ , and due to cyclic and linear properties of tr function (Fact 4). Let we denote the matrix  $X X^T$  by  $P$ , and  $(A^T A - \tilde{A}^T \tilde{A})$  by matrix  $M$ . Thus, the problem reduces to bounding the term  $\text{tr}(PM)$ . Let  $\lambda_i(M)$  is the  $i$ th eigenvalue, and  $\{w_i\}_{i=1}^d$  be the eigenvectors of  $M$ , then  $M = \sum_{i=1}^d \lambda_i(M) w_i w_i^T$ . The following expression holds due to linearity of trace function.

$$\text{tr}(PM) = \text{tr}(P \sum_{i=1}^d \lambda_i(M) w_i w_i^T) = \sum_{i=1}^d \lambda_i(M) \text{tr}(P w_i w_i^T)$$

Further, we bound the summation  $\sum_{i=1}^d \text{tr}(P w_i w_i^T)$ ,

$$\begin{aligned} &\sum_{i=1}^d \text{tr}(P w_i w_i^T) \\ &= \text{tr}(P W W^T) = \text{tr}(P^T P W W^T W W^T) \\ &= \text{tr}(P^T P W W^T) = \text{tr}(P W W^T P^T) \\ &= \|PW\|_F^2 \leq \|P\|_F^2 = \|X X^T\|_F^2 \leq j \end{aligned} \quad (7)$$

where,  $W \in \mathbb{R}^{d \times d}$  having columns as eigenvectors of  $M$ . The above equalities follow as  $P = X X^T$ , then  $P^T P = P$ ; similarly  $W W^T = W W^T W W^T$ , also  $X^T X = I, W^T W = I$ . Finally, the inequality  $\|PW\|_F^2 \leq \|P\|_F^2$  follows from Fact 7.

Further,  $P = X X^T$  has all singular values either 1 or 0.

$$0 \leq \text{tr}(P w_i w_i^T) = w_i^T P w_i \leq \|w_i\|_2^2 \|P\|_2^2 \leq 1 \quad (8)$$

Thus, for  $1 \leq i \leq d$ ,  $\text{tr}(P w_i w_i^T)$  has  $d$  values, and each value is at most 1 (Equation 8), and sum of all of them is at most  $j$  (Equation 7).

Further as  $M$  is symmetric, we have

$$M = A^T A - \tilde{A}^T \tilde{A} = \left( A - A R' R'^T \right)^T \left( A - A R' R'^T \right) \quad (9)$$

Equality 9 holds due to Fact 3. Equality 9 and Fact 5 shows that  $M$  is a positive semidefinite matrix, and as a consequence it has all nonnegative eigenvalues. Then, the summation  $\sum_{i=1}^d \lambda_i(M) \text{tr}(P w_i w_i^T)$  is maximized when  $\text{tr}(P w_i w_i^T) = 1$ , that is for those eigenvectors which corresponds to largest magnitude eigenvalues of  $M$ . Thus,

$$0 \leq \text{tr}(PM) = \sum_{i=1}^d \lambda_i(M) \text{tr}(P w_i w_i^T) \leq \sum_{i=1}^j \lambda_i(M).$$

As a consequence, we have  $0 \leq \text{tr}(XX^T(A^T A - \tilde{A}^T \tilde{A})) \leq \sum_{i=1}^j \lambda_i(A^T A - \tilde{A}^T \tilde{A}) = \sum_{i=1}^j \sigma_i^2(A - \tilde{A}) = \|(A - \tilde{A})_j\|_F^2$ . Here matrix  $(A - \tilde{A})_j$  is the matrix restricted to rank  $j$  of the matrix  $A - \tilde{A}$ . Further, as matrix  $\tilde{A}$  is of rank at most  $m$ ,  $\tilde{A} + (A - \tilde{A})_j$  is of rank at most  $m + j$ . Thus, we have

$$\|A - (\tilde{A} + (A - \tilde{A})_j)\|_F^2 \geq \sum_{i=m+j+1}^d \sigma_i^2 \quad (10)$$

$$\|A - \tilde{A}\|_F^2 - \|(A - \tilde{A})_j\|_F^2 \geq \sum_{i=m+j+1}^d \sigma_i^2 \quad (11)$$

$$\begin{aligned} \|(A - \tilde{A})_j\|_F^2 &\leq \|A - \tilde{A}\|_F^2 - \sum_{i=m+j+1}^d \sigma_i^2 \\ \|(A - \tilde{A})_j\|_F^2 &\leq (1 + \epsilon) \|A - A^{(m)}\|_F^2 - \sum_{i=m+j+1}^d \sigma_i^2 \end{aligned} \quad (12)$$

$$\begin{aligned} &= (1 + \epsilon) \sum_{i=m+1}^d \sigma_i^2 - \sum_{i=m+j+1}^d \sigma_i^2 \\ &= \sum_{i=m+1}^{m+j} \sigma_i^2 + \epsilon \sum_{i=m+1}^d \sigma_i^2 \\ &\leq \epsilon \sum_{j+1}^d \sigma_i^2 + \epsilon \sum_{i=j+1}^d \sigma_i^2 \\ &= 2\epsilon \|AX^\perp\|_F^2 \end{aligned} \quad (13)$$

Inequality 10 follows from Fact 1; Inequality 11 follows from Fact 3, where  $(A - \tilde{A})_j = (A - \tilde{A})XX^T$  for an orthonormal matrix  $X \in \mathbb{R}^{d \times j}$ ; Inequality 12 follows from Theorem 2; Inequality 13 holds as  $m = \lceil \frac{j}{\epsilon} \rceil$ , and due to the following:

$$\sum_{i=m+1}^{m+j} \sigma_i^2 \leq j\sigma_{m+1}^2 = \epsilon m \sigma_{m+1}^2 \leq \epsilon m \sigma_{j+1}^2 \leq \epsilon \sum_{i=j+1}^{m+j} \sigma_i^2 \leq \epsilon \sum_{i=j+1}^d \sigma_i^2.$$

## 4 Faster coresnet construction for projective clustering

In this section, extending the result (Theorem 3) of the previous section, we present a randomized coresnet construction for the problem of projective clustering. More precisely, if  $L_1, \dots, L_k$  be a set of  $k$  subspaces each of dimension at most  $j$ , and let  $\mathcal{C}$  be a closed set containing union of them, then our randomized coresnet is a matrix of very small rank (independent of  $d$ ) and it approximately preserve the distances from every such closed set  $\mathcal{C}$ , with high probability. Our main contribution is the dimensionality reduction step of the coresnet construction which is presented in Algorithm 1.

**Proof of Theorem 1:** Let  $[X^*]_{d \times j^*}$  be a matrix with orthonormal columns whose span is  $L^*$ , and let  $L^{*\perp}$  be the orthogonal complement of  $L^*$  spanned by  $[X^{*\perp}]_{d \times (d-j^*)}$ . In our analysis, we use the equality,  $\text{dist}^2(A, \mathcal{C}) = \|AX^{*\perp}\|_F^2 + \text{dist}^2(AX^*X^{*T}, \mathcal{C})$ , which holds true from the Pythagorean theorem. We have

$$\begin{aligned} &|(\text{dist}^2(A^*, \mathcal{C}) + \Delta^*) - \text{dist}^2(A, \mathcal{C})| \\ &= \left| \left( \|A^*X^{*\perp}\|_F^2 + \text{dist}^2(A^*X^*X^{*T}, \mathcal{C}) + \Delta^* \right) - \left( \|AX^{*\perp}\|_F^2 + \text{dist}^2(AX^*X^{*T}, \mathcal{C}) \right) \right| \\ &\leq \underbrace{\left| \left( \|A^*X^{*\perp}\|_F^2 + \Delta^* - \|AX^{*\perp}\|_F^2 \right) \right|}_{\text{first term}} + \underbrace{\left| \left( \text{dist}^2(A^*X^*X^{*T}, \mathcal{C}) - \text{dist}^2(AX^*X^{*T}, \mathcal{C}) \right) \right|}_{\text{second term}} \end{aligned}$$

- 1 **Input:**  $A \in \mathbb{R}^{n \times d}$ , an integer  $1 \leq j < d - 1$ , and an integer  $1 \leq k < n - 1$  such that  $j^* \leq d - 1$ , where  $j^* = k(j + 1)$ ,  $\epsilon \in (0, 1)$ ,  $\delta \in (0, 1)$ .
- 2 **Result:** Dimensionality reduction for randomized coresets construction for the projective clustering.
- 3 Compute an Johnson-Lindenstrauss matrix  $[\mathcal{S}]_{r \times n}$  having *i.i.d.*  $\pm 1$  entries and zero-mean, where  $r = O((\frac{m^*}{\epsilon} + m^* \log m^*) \log \frac{1}{\delta})$ ,  $m^* = \lceil \frac{52j^*}{\epsilon^2} \rceil$ .
- 4 Compute the matrix  $\pi_{SA}(A)$ .
- 5 Compute the SVD of  $\pi_{SA}(A)$ , let  $R^* \in \mathbb{R}^{d \times m^*}$  be the first  $m^*$  right singular vectors of  $\pi_{SA}(A)$ .
- 6 Let us denote  $AR^*R^{*T}$  by  $A^*$ , and output  $A^*$ .

**Algorithm 1:** Algorithm for dimensionality reduction for projective clustering.

We have two terms to bound in the above expression. The first term can be upper bounded using a similar analysis as of Theorem 3 which holds true with probability at least  $1 - \delta$ . (In Theorem 3, we replace  $j$  by  $j^*$ ,  $m$  by  $m^*$ ,  $\epsilon$  by  $\frac{\epsilon^2}{52}$ , and  $\Delta$  by  $\Delta^*$ .)

$$\left| \|A^* X^{*\perp}\|_F^2 + \Delta^* - \|AX^{*\perp}\|_F^2 \right| \leq \frac{\epsilon^2}{26} \|AX^{*\perp}\|_F^2 \quad (14)$$

In order to bound the second term  $\left| \text{dist}^2(A^* X^* X^{*T}, \mathcal{C}) - \text{dist}^2(AX^* X^{*T}, \mathcal{C}) \right|$ , we use a triangle inequality from Lemma 10. For any  $\varepsilon \in (0, 1)$  and from Lemma 10, we have

$$\begin{aligned} & \left| \text{dist}^2(A^* X^* X^{*T}, \mathcal{C}) - \text{dist}^2(AX^* X^{*T}, \mathcal{C}) \right| \\ & \leq \frac{12}{\varepsilon} \|A^* X^* X^{*T} - AX^* X^{*T}\|_F^2 + \frac{\varepsilon}{2} \text{dist}^2(AX^* X^{*T}, \mathcal{C}) \\ & \leq \frac{12}{\varepsilon} \left( \frac{\epsilon^2}{26} \|AX^{*\perp}\|_F^2 \right) + \frac{\varepsilon}{2} \text{dist}^2(AX^* X^{*T}, \mathcal{C}) \\ & \leq \frac{12}{\varepsilon} \left( \frac{\epsilon^2}{26} \|AX^{*\perp}\|_F^2 \right) + \frac{\varepsilon}{2} \text{dist}^2(A, \mathcal{C}) \end{aligned} \quad (15)$$

Inequality 15 holds due to Lemma 12. Thus, we have

$$\begin{aligned} & \left| \text{dist}^2(A^* X^* X^{*T}, \mathcal{C}) - \text{dist}^2(AX^* X^{*T}, \mathcal{C}) \right| \\ & \leq \frac{12}{\varepsilon} \left( \frac{\epsilon^2}{26} \|AX^{*\perp}\|_F^2 \right) + \frac{\varepsilon}{2} \text{dist}^2(A, \mathcal{C}) \end{aligned} \quad (16)$$

Equation 14, in conjunction with Equation 16, gives us the following:

$$\begin{aligned} & |(\text{dist}^2(A^*, \mathcal{C}) + \Delta^*) - \text{dist}^2(A, \mathcal{C})| \\ & \leq \left( 1 + \frac{12}{\varepsilon} \right) \frac{\epsilon^2}{26} \|AX^{*\perp}\|_F^2 + \frac{\varepsilon}{2} \text{dist}^2(A, \mathcal{C}) \\ & \leq \left( 1 + \frac{12}{\varepsilon} \right) \frac{\epsilon^2}{26} \text{dist}^2(A, \mathcal{C}) + \frac{\varepsilon}{2} \text{dist}^2(A, \mathcal{C}) \end{aligned}$$

$$\begin{aligned}
&= \left( \frac{\epsilon^2}{26} + \frac{12\epsilon^2}{26\epsilon} + \frac{\epsilon}{2} \right) \text{dist}^2(A, \mathcal{C}) \\
&= \left( \frac{\epsilon^2}{26} + \frac{12\epsilon}{26} + \frac{\epsilon}{2} \right) \text{dist}^2(A, \mathcal{C}) \\
&\leq \epsilon \text{dist}^2(A, \mathcal{C})
\end{aligned} \tag{17}$$

Equality 17 holds by choosing  $\epsilon = \epsilon$ , and as  $\epsilon^2/26 + 12\epsilon/26 < \epsilon/2$ .

**Lemma 12** *Let  $X^* \in \mathbb{R}^{d \times j^*}$  be a matrix with orthonormal columns whose span is  $L^*$ , then in Algorithm 1 the following is true with probability at least  $1 - \delta$*

$$\|A^* X^* X^{*T} - A X^* X^{*T}\|_F^2 \leq \frac{\epsilon^2}{26} \|A X^{*\perp}\|_F^2.$$

*Proof.*

$$\begin{aligned}
\|A^* X^* X^{*T} - A X^* X^{*T}\|_F^2 &= \|A X^* X^{*T} - A^* X^* X^{*T}\|_F^2 \\
&= \|(A - A^*) X^* X^{*T}\|_F^2 \\
&= \text{tr} \left( (A - A^*) X^* X^{*T} \left( (A - A^*) X^* X^{*T} \right)^T \right) \\
&= \text{tr} \left( (A - A^*) X^* X^{*T} X^* X^{*T} (A - A^*)^T \right) \\
&= \text{tr} \left( (A - A^*)^T (A - A^*) X^* X^{*T} \right) \tag{18}
\end{aligned}$$

$$= \text{tr} \left( X^* X^{*T} (A^T A - A^{*T} A^*) \right) \tag{19}$$

$$\leq \frac{\epsilon^2}{26} \|A X^{*\perp}\|_F^2 \tag{20}$$

Equality 18 holds due to cyclic property of  $\text{tr}$  function and  $X^{*T} X^* = I$ ; Equality 19 holds due to cyclic property of  $\text{tr}$ , and  $(A - A^*)^T (A - A^*) = (A - A R^* R^{*T})^T (A - A R^* R^{*T}) = A^T A - A^{*T} A^*$  (after simplification, see Fact 3); finally Inequality 20 holds by following the steps of proof of Lemma 11 (from Equation 6), where we replace  $\epsilon$  by  $\frac{\epsilon^2}{52}$ ;  $j$  by  $j^*$ ,  $A$  by  $A^*$ ,  $X$  by  $X^*$ .

*Remark 3.* Please note that it suffices to store the matrix  $A R^*$  which is of dimension  $m^*$ , where  $m^* = O\left(\frac{k(j+1)}{\epsilon^2}\right)$ . However, for the purpose of our analysis, we use the matrix  $A R^* R^{*T}$  which is of dimension  $d$ , and rank  $m^*$ . Further, the space required to store our coreset is  $O(nm^* + 1)$ - we need  $O(nm^*)$  to store the matrix  $A R^*$ , and  $O(1)$  space to store the term  $\Delta^*$ ; on the other hand, the space required to store  $A$  is  $O(nd)$ .

**Comparison with coreset construction of [10]:** Coreset construction of [10] require projecting the rows of  $A$  on its first  $O(k(j+1)/\epsilon^2)$  right singular vectors which gives a matrix of rank  $O(k(j+1)/\epsilon^2)$  and it approximately preserve the distance from any closed  $\mathcal{C}$ . Their construction required computing SVD of

the given matrix  $A$ , whose complexity is  $\min\{n^2d, nd^2\}$ . In our construction, we showed that it is also suffices to project the rows of  $A$  on  $O(k(j+1)/\epsilon^2)$  orthonormal vectors that closely approximate the right singular vectors of  $A$ . We now give an expected time bound on the running time of Algorithm 1. Time required for execution of line number 3, 4, 5 is

$$\begin{aligned} & O\left(\mathbf{nnz}(A)\left(\frac{m^*}{\epsilon} + m^* \log m^*\right) + (n+d)\left(\frac{m^*}{\epsilon} + m^* \log m^*\right)^2\right) \\ &= O\left(\mathbf{nnz}(A)\left(\frac{j^*}{\epsilon^3} + \frac{j^*}{\epsilon^2} \log \frac{j^*}{\epsilon^2}\right) + (n+d)\left(\frac{j^*}{\epsilon^3} + \frac{j^*}{\epsilon^2} \log \frac{j^*}{\epsilon^2}\right)^2\right), \end{aligned}$$

due to [19], where  $j^* = k(j+1)$ . Further, line number 6 requires time - for projecting  $A$  on  $R^*$ , which due to an elementary matrix multiplication is  $O(ndm^*) = O\left(\frac{ndj^*}{\epsilon^2}\right)$ . Thus, total expected running time of Algorithm 1 is

$$\begin{aligned} & O(\mathbf{nnz}(A)\left(\frac{j^*}{\epsilon^3} + \frac{j^*}{\epsilon^2} \log \frac{j^*}{\epsilon^2}\right) + (n+d)\left(\frac{j^*}{\epsilon^3} + \frac{j^*}{\epsilon^2} \log \frac{j^*}{\epsilon^2}\right)^2 + \frac{ndj^*}{\epsilon^2}) \\ &= \tilde{O}\left(\mathbf{nnz}(A)\frac{j^*}{\epsilon^3} + (n+d)\frac{j^{*2}}{\epsilon^6} + \frac{ndj^*}{\epsilon^2}\right). \end{aligned}$$

Clearly, if  $n \geq d$  and  $j^* = o(n)$ , or, if  $n < d$  and  $j^* = o(d)$ , then our expected running time is better than that of [10].

As a corollary of Theorem 1, and using the known techniques from [10,20,8] on  $A^*$ , we present the cardinality reduction step of cores et construction as follows:

**Corollary 13 (Corollary 9.1 of [10])** *Let  $A \in \{1, 2, \dots, \Lambda\}^{n \times d}$ , with  $\Lambda \in (nd)^{O(1)}$ ,  $d \in n^{O(1)}$ . There is a matrix  $\mathcal{Q} \in \mathbb{R}^{l \times d'}$  with  $l = \text{poly}(2^{kj}, \frac{1}{\epsilon}, \log n, \log \Lambda)$ ,  $d' = O(k(j+1)/\epsilon^2)$ ; and a weight function associated with the rows of  $\mathcal{Q}$ , i.e.  $w : \mathcal{Q}_{i^*} \rightarrow [0, \infty)$  such that for every closed set  $\mathcal{C}$ , which is the union of  $k$  affine  $j$ -subspaces of  $\mathbb{R}^d$ , the following holds with high probability*

$$(1 - \epsilon)\Sigma_{i=1}^n \text{dist}^2(A_{i^*}, \mathcal{C}) \leq \Sigma_{i=1}^l w(\mathcal{Q}_{i^*}) \text{dist}^2(\mathcal{Q}_{i^*}, \mathcal{C}) \leq (1 + \epsilon)\Sigma_{i=1}^n \text{dist}^2(A_{i^*}, \mathcal{C}).$$

As a corollary of Theorem 1, we obtain following randomized cores et result for  $k$ -mean clustering.

**Corollary 14** *Let  $A \in \mathbb{R}^{n \times d}$ ,  $\epsilon \in (0, 1)$ , and  $k$  an integer less than  $(d-1)$  and  $(n-1)$ . Then there is a randomized algorithm which outputs a matrix  $A'$  of rank  $O(k/\epsilon^2)$  such that for every set of  $k$  points  $\{c_i\}_{i=1}^k \in \mathbb{R}^d$  represented as the rows of matrix  $C$ , the following holds with high probability:*

$$|(\text{dist}^2(A', C) + \Delta') - \text{dist}^2(A, C)| \leq \epsilon \text{dist}^2(A, C).$$

*The expected running time of the algorithm is  $\tilde{O}\left(\mathbf{nnz}(A)\frac{k}{\epsilon^3} + (n+d)\frac{k^2}{\epsilon^6} + \frac{ndk}{\epsilon^2}\right)$ .*

*Where,  $m' = O(k/\epsilon^2)$ ;  $\Delta' = \|A - A^{(m')}\|_F^2$ ; and  $\text{dist}^2(A, C)$  denotes the sum of square distances from each row of  $A$  to its closest point in  $C$ .*

## 5 Conclusion and open problems

We presented a randomized coresets construction for projective clustering *via* low rank approximation. We first revisit the result of [5] for the subspace clustering, and then extend their result to construct a randomized coresets for projective clustering. We showed that our construction is significantly faster (when the values of  $k$  and  $j$  are small), as compared to the corresponding deterministic construction of [10], and it also maintains nearly the same accuracy. Our work leaves several open problems - improving dimensionality reduction bounds for projective clustering, or giving a matching lower bound for the same.

## References

1. P. K. Agarwal, S. Har-Peled, and K. R. Varadarajan. Approximating extent measures of points. *J. ACM*, 51(4):606–635, 2004.
2. P. K. Agarwal, S. Har-Peled, and K. R. Varadarajan. geometric approximation via coresets. *Current Trends in Combinatorial and Computational Geometry (E. Welzl, ed.)*, 2007.
3. C. C. Aggarwal, C. M. Procopiuc, J. L. Wolf, P. S. Yu, and J. S. Park. Fast algorithms for projected clustering. In *SIGMOD 1999, Proceedings ACM SIGMOD International Conference on Management of Data, June 1-3, 1999, Philadelphia, Pennsylvania, USA.*, pages 61–72, 1999.
4. K. Chakrabarti and S. Mehrotra. Local dimensionality reduction: A new approach to indexing high dimensional spaces. In *VLDB 2000, Proceedings of 26th International Conference on Very Large Data Bases, September 10-14, 2000, Cairo, Egypt*, pages 89–100, 2000.
5. M. B. Cohen, S. Elder, C. Musco, C. Musco, and M. Persu. Dimensionality reduction for k-means clustering and low rank approximation. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015*, pages 163–172, 2015.
6. M. Ester, H. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96), Portland, Oregon, USA*, pages 226–231, 1996.
7. D. Feldman, A. Fiat, and M. Sharir. Coresets for weighted facilities and their applications. In *47th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2006), 21-24 October 2006, Berkeley, California, USA, Proceedings*, pages 315–324, 2006.
8. D. Feldman and M. Langberg. A unified framework for approximating and clustering data. In *Proceedings of the 43rd ACM Symposium on Theory of Computing, STOC 2011, San Jose, CA, USA, 6-8 June 2011*, pages 569–578, 2011.
9. D. Feldman, M. Monemizadeh, C. Sohler, and D. P. Woodruff. Coresets and sketches for high dimensional subspace approximation problems. In *Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms*, pages 630–649, 2010.
10. D. Feldman, M. Schmidt, and C. Sohler. Turning big data into tiny data: Constant-size coresets for k-means, pca and projective clustering. In *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1434–1453, 2013.

11. S. Har-Peled. No, coreset, no cry. In *FSTTCS 2004: Foundations of Software Technology and Theoretical Computer Science, 24th International Conference, Chennai, India, December 16-18, 2004, Proceedings*, pages 324–335, 2004.
12. S. Har-Peled and S. Mazumdar. On coresets for k-means and k-median clustering. In *Proceedings of the 36th Annual ACM Symposium on Theory of Computing, Chicago, IL, USA, June 13-16, 2004*, pages 291–300, 2004.
13. M. R. Henzinger, P. Raghavan, and S. Rajagopalan. Computing on data streams. In *External Memory Algorithms, Proceedings of a DIMACS Workshop, New Brunswick, New Jersey, USA, May 20-22, 1998*, pages 107–118, 1998.
14. A. Hinneburg, C. C. Aggarwal, and D. A. Keim. What is the nearest neighbor in high dimensional spaces? In *VLDB 2000, Proceedings of 26th International Conference on Very Large Data Bases, September 10-14, 2000, Cairo, Egypt*, pages 506–515, 2000.
15. W. B. Johnson and J. Lindenstrauss. Extensions of lipschitz mappings into a hilbert space. *Conference in modern analysis and probability (New Haven, Conn., 1982)*, Amer. Math. Soc., Providence, R.I., pages 189–206, 1983.
16. S. Lloyd. Least squares quantization in pcm. *IEEE Trans. Inf. Theor.*, 28(2):129–137, Sept. 2006.
17. M. W. Mahoney. Randomized algorithms for matrices and data. *Foundations and Trends in Machine Learning*, 3(2):123–224, 2011.
18. J. M. Phillips. Coresets and sketches. *CoRR*, abs/1601.00617, 2016.
19. T. Sarlós. Improved approximation algorithms for large matrices via random projections. In *47th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2006), 21-24 October 2006, Berkeley, California, USA, Proceedings*, pages 143–152, 2006.
20. K. R. Varadarajan and X. Xiao. A near-linear algorithm for projective clustering integer points. In *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2012, Kyoto, Japan, January 17-19, 2012*, pages 1329–1342, 2012.
21. T. Zhang, R. Ramakrishnan, and M. Livny. BIRCH: A new data clustering algorithm and its applications. *Data Min. Knowl. Discov.*, 1(2):141–182, 1997.